2	In this appendix, we provide details skipped in the main text. The content is organized as follows:
3 4	• Section 1. Detailed algorithms of GPS, including binary search and simulation process. ( <i>c.f.</i> §5.4 of the main text)
5	• Section 2. Validation of the switching point existence. (c.f. §4.3 of the main text)
6 7	• Section 3. Validation on the properties of the simulation methods. ( <i>c.f.</i> §5.2 of the main text)
8 9	• Section 4. Additional experimental results and ablation studies. ( <i>c.f.</i> §6.2, §6.5 and §6.6 of the main text)
10 11	• Section 5. Exploration of new base policies based on curriculum learning, and their performance with GPS.
12	• Section 6. Experimental details and hyperparameters. ( <i>c.f.</i> §6.1 of the main text)

Appendix

# **13 1 Detailed Algorithms**

#### 14 **1.1 Simulation Process**

1

The pseudocode of our simulation process (Fig. 2 in the main text) is listed in Algorithm 1. We use the notation  $P_{1:i}$  to represent the task distributions from  $P_1$  to  $P_i$ . Likewise, we use the notation  $s_{1:i}$  to represent the switching point from  $s_1$  to  $s_i$ . The memory construction function BuildM takes as arguments the previous memory and a list of switching points. This function internally checks whether the list of switching points are enough to construct the current memory. If yes, it constructs the memory as described in §4; otherwise, it utilizes the described pseudo-memory construction methods as described in §5.

#### Algorithm 1 GlobalSim

**Input:** Tested point  $a_i$ ; Pseudo-task distributions  $\tilde{P}_{(i+1):T}$ ; Switching points  $s_{1:(i-1)}$ ; Local updating method g; Current model parameters  $\theta_i$  and current memory  $\mathcal{M}_{i-1}$ . Initialize  $\tilde{\theta}_{i:i} \leftarrow \theta_i$ Initialize memory  $\tilde{\mathcal{M}}_i \leftarrow \text{BuildM}(\mathcal{M}_{i-1}, s_{1:(i-1)} \cup \{a_i\})$   $j \leftarrow i + 1$ while  $j \leq T$  do Local update:  $\tilde{\theta}_{i:j} \leftarrow g(\tilde{\theta}_{i:(j-1)}, \tilde{P}_j, \tilde{\mathcal{M}}_{j-1})$ Build memory:  $\tilde{\mathcal{M}}_j \leftarrow \text{BuildM}(\tilde{\mathcal{M}}_{j-1}, s_{1:i})$   $j \leftarrow j + 1$ end while Compute loss:  $l \leftarrow \mathbb{E}_{(x_j, y_j) \sim P_j} \ell(y_j, f(x_j; \tilde{\theta}_{j:T}))$ return Loss: l

#### 22 1.2 Binary Search

In the global binary search as listed in Algorithm 2, we take a search stride  $\epsilon = 20$  to increase the robustness of the algorithm.

# 25 **1.3 GPS Algorithm**

<sup>26</sup> Based on the simulation process and binary search, we describe our Global Pseudo-task Simulation

<sup>27</sup> method in Algorithm 3.

Algorithm 2 GlobalBS

**Input:** Number of tasks T; Task distributions  $P_{1:i}$ ; Switching points  $s_{1:(i-1)}$ ; Local updating method g; Current model parameters  $\theta_i$  and current memory  $\mathcal{M}_{i-1}$ ; Search stride  $\epsilon$ . Synthesize pseudo-tasks from  $P_i$  with task distributions.  $start \leftarrow 0$ end  $\leftarrow |\mathcal{M}|/i$ Loss dictionary:  $loss\_dict \leftarrow \varnothing$ while  $end - start \ge \epsilon \, \mathbf{do}$  $next \leftarrow (start + end)/2$ if *next* not in *loss\_dict* then  $loss \leftarrow \text{GlobalSim}(next, \tilde{P}_{(i+1):T}, s_{1:(i-1)}, g, \theta_i, \mathcal{M}_{i-1})$  $loss\_dict \leftarrow loss\_dict \cup \{next: loss\}$ else  $loss \leftarrow loss\_dict[next]$ end if if  $next - \epsilon$  not in  $loss\_dict$  then  $left_loss \leftarrow \text{GlobalSim}(next - \epsilon, P_{(i+1):T}, s_{1:(i-1)}, g, \theta_i, \mathcal{M}_{i-1})$  $loss\_dict \leftarrow loss\_dict \cup \{next - \epsilon : left\_loss\}$ else  $left\_loss \leftarrow loss\_dict[next - \epsilon]$ end if if  $next + \epsilon$  not in  $loss\_dict$  then  $right_loss \leftarrow GlobalSim(next + \epsilon, \tilde{P}_{(i+1):T}, s_{1:(i-1)}, g, \theta_i, \mathcal{M}_{i-1})$  $loss\_dict \leftarrow loss\_dict \cup \{next + \epsilon : right\_loss\}$ else  $right_loss \leftarrow loss_dict[next + \epsilon]$ end if if  $left_loss < loss$  then  $end \gets next$ continue else if *right\_loss < loss* then  $start \leftarrow next$ continue else  $s_i \leftarrow next$ break end if end while if  $s_i$  is not assigned then  $s_i \leftarrow argmin_{loss}(loss\_dict)$ end if return Switching point:  $s_i$ 

Algorithm 3 Global Pseudo-task Simulation (GPS)

**Input:** Number of tasks *T*; Task distributions  $P_i, i \in \mathcal{T}$ ; Local updating method  $g(\cdot)$ . Initialize parameters  $\theta_0$ Initialize memory  $\mathcal{M}_0 = \emptyset$   $i \leftarrow 1$ while  $i \leq T$  do Local update:  $\theta_i \leftarrow g(\theta_{i-1}, P_i, \mathcal{M}_{i-1})$ Find switching point:  $s_i \leftarrow \text{GlobalBS}(T, P_{1:i}, s_{1:(i-1)}, g, \theta_i, \mathcal{M}_{i-1})$ Build memory:  $\mathcal{M}_i \leftarrow \text{BuildM}(\mathcal{M}_{i-1}, s_{1:i})$   $i \leftarrow i + 1$ end while return Model parameters:  $\theta_T$ 

# 28 2 Validation of the Switching Point Existence

As a validation of our switching point existence, we further show the switching point of other benchmarks in our experiments. For each benchmark, we plot the global loss  $\mathcal{L}_G$  as a function of  $a_i$ and select 5 different tasks  $t_i$ . Each plot shows clearly the switching point in Fig. 1.



Figure 1: Switching points of the first 5 tasks in three evaluation benchmarks: TinyImageNet, S-CIFAR-100, P-MNIST. We show the change of the global loss,  $\mathcal{L}_G$  w.r.t. the ratio of ER-Ring-Full in the memory.

31

# 32 **3** Validation of the Simulation Method

In this section, we provide the empirical supporting evidences for our hypotheses of the simulation
 method.

#### 35 3.1 Task Difficulties

- <sup>36</sup> First, we show the task difficulties in the evaluated benchmarks have small variations, as in Table 1.
- <sup>37</sup> For P-MNIST, S-CIFAR-100 and TinyImageNet, we evaluate the first 5 tasks end-to-end for simplicity. For S-CIFAR-10, we evaluate all the tasks end-to-end.

Table 1: Accuracy and variance of accuracy of tasks from four vision benchmarks trained end-to-end.

Dataset	Task 1	Task 2	Task 3	Task 4	Task 5	Variance
P-MNIST	97.48	97.28	97.33	97.78	97.53	0.03
S-CIFAR-10	98.20	94.85	96.50	98.90	98.15	2.18
S-CIFAR-100	85.70	87.70	88.10	88.10	86.20	1.02
TinyImageNet	78.20	76.80	77.30	76.50	77.20	0.33

38

<sup>39</sup> Further, we evaluate the difficulty along the pseudo-task sequence synthesized from the first task of

40 S-CIFAR-100 by permutation, rotation and blurring, we generate 5 tasks for each simulation method.

41 The results in Table 2 shows permutation and rotation generate tasks with similar difficulties as the

42 original S-CIFAR-100 tasks, while blurring generate tasks with increasing difficulties as the task

43 sequence grows.

Method	Task 1	Task 2	Task 3	Task 4	Task 5	Variance
Permutation Rotation	85.50 85.40	85.70 85.40	85.30 85.70	86.10 84.90	85.10 84.50	0.11 0.18
Blurring	83.90	81.70	78.90	74.40	69.50	26.80

Table 2: Accuracy and variance of accuracy of pseudo-tasks synthesized by different methods from the first task of S-CIFAR-100.

#### 44 3.2 Forward Transfer Ability of Tasks

<sup>45</sup> Next, we explore the forward transfer ability of the pseudo-tasks *vs*. the real tasks. We evaluate <sup>46</sup> the model trained on the first two tasks of S-CIFAR-100 on task sequences different by different

47 simulation methods.

Table 3 shows that the permutation pseudo-tasks and the real tasks both allows zero transfer ability, as the random guess accuracy of a 10-class classification is 10%. Rotation, instead, creates a task

sequence that allows nearly perfect forward transfer ability. Blurring creates a task sequence which

allows some forward transfer ability from the beginning, but it gradually reduces to a random guess as task difficulty grows.

Table 3: Forward transfer ability of different simulation methods after training task  $t_1$  and task  $t_2$  on S-CIFAR-100. Numbers are the accuracy of the task.

Dataset	Method	Task 3	Task 4	Task 5	Task 6	Task 7
S-CIFAR-100	Real	11.80	10.40	10.30	9.40	9.70
	Permutation	10.50	10.40	10.30	10.10	11.10
	Rotation	75.40	78.10	77.70	79.90	80.50
	Blurring	70.90	65.70	52.90	30.40	15.50

52

#### 53 3.3 Global Loss w.r.t. Single Task Loss

<sup>54</sup> The empirical analysis on four benchmarks suggests the summed global loss and the single final loss

of a task  $t_j$  are positively correlated as a function of  $a_j$ . Specifically, the switching point for both functions are similar. Fig. 2 shows four tasks in four benchmarks, respectively. This observation

<sup>57</sup> implies our simulation objective well approximate the global objective as in the offline setting.



Figure 2: Global loss w.r.t. single task loss with four tasks on four benchmarks.

## 58 4 Additional Experimental Results

#### 59 4.1 Accuracy of GPS with Small Memory Buffer

We have also carried out experiments of GPS to evaluate its performance when the memory buffer is relatively small, as shown in Table 4. With a small size memory buffer, GPS does not show significant improvement. One reason is the switching point  $s_i$  is very close to 1, *i.e.*, taking the pure

- 63 ER-Ring-Full policy is good enough. Another cause is our base policy assumption does not work
- 64 very well under the small memory buffer size as model becomes more sensitive to points selection

65 under small  $\mathcal{M}$ .

Table 4: Accuracy of GPS using permutation and ER baselines on four datasets with smaller buffer sizes.

Method   <i>M</i>	<b>P-MNIST</b> 100	<b>S-CIFAR-10</b> 20	S-CIFAR-100 200	<b>TinyImageNet</b> 200
ER-Res ER-Ring-Full	$65.59_{\pm 1.38}$ $66.10_{\pm 1.36}$	$\begin{array}{c} 80.68 \scriptstyle{\pm 2.28} \\ 81.30 \scriptstyle{\pm 1.98} \end{array}$	$\begin{array}{c} 64.99{\scriptstyle\pm1.74} \\ 65.95{\scriptstyle\pm0.96} \\ \end{array}$	$38.60{\scriptstyle\pm0.74}\ 40.85{\scriptstyle\pm0.78}$
ER-Hybrid GPS	$\begin{array}{c} 66.30 \scriptstyle \pm 1.21 \\ 66.50 \scriptstyle \pm 1.11 \end{array}$	$\frac{81.43{\scriptstyle\pm2.54}}{81.78{\scriptstyle\pm1.56}}$	$66.30{\scriptstyle\pm1.32} \\ 66.51{\scriptstyle\pm0.88}$	$\begin{array}{c} 40.75 {\scriptstyle \pm 0.54} \\ 40.89 {\scriptstyle \pm 0.45} \end{array}$

#### 66 4.2 Results of GPS with DER++, HAL and Baselines

<sup>67</sup> We put complete results of GPS+DER, GPS+DER++, GPS+HAL and baselines in Table 5. Note the

<sup>68</sup> A-GEM [5], iCaRL [9] and GSS [1] use the same memory size as the ER series for fair comparison.

<sup>69</sup> The results stand as a complete empirical support to illustrate that the performance of other ER

variants have been improved after using our GPS method. We also reported the results of GPS+DER,

71 GPS+DER++ and GPS+HAL with smaller memory sizes in Table 6.

Table 5: Accuracy of GPS using permutation incorporating DER, DER++ [3] and HAL [4], comparing to other methods. '-' indicates experiments we were unable to run, due to compatibility issues (e.g. Domain IL for iCaRL) or intractable training time or memory utilization (e.g. OGD, GSS on TinyImageNet).

	P-MNIST	S-CIFAR10	S-CIFAR100	TinyImageNet
oEWC	$69.21_{\pm 2.92}$	$62.97 \pm 3.55$	$55.37_{\pm 2.71}$	$20.81 \pm 0.95$
iCaRL	-	$88.97 \pm 2.77$	$78.21 \pm 1.01$	$38.77 {\scriptstyle \pm 3.68}$
GSS	$86.34 \pm 4.28$	$87.80 {\pm 2.71}$	$77.34_{\pm 3.21}$	-
A-GEM	$77.36 \pm 1.28$	$83.87 \pm 1.55$	$69.61 \pm 1.47$	$25.30 \pm 0.87$
OGD	$81.52_{\pm 2.21}$	-	-	-
	P-MNIST	S-CIFAR10	S-CIFAR100	TinyImageNet
$ \mathcal{M} $	1000	200	2000	2000
HAL	$87.69 \pm 0.34$	$89.29 \pm 1.31$	$80.81 \pm 1.21$	$61.27 \pm 1.10$
GPS+HAL	$88.73 \scriptstyle \pm 0.03$	$91.27 \scriptstyle \pm 0.93$	$82.33{\scriptstyle \pm 0.47}$	$63.24 \scriptstyle \pm 0.80$
DER	$90.47 \pm 2.69$	$91.04 \pm 0.18$	$81.78 \pm 0.50$	$60.90 \pm 1.08$
GPS+DER	$90.27_{\pm 1.78}$	$91.53_{\pm 0.13}$	$83.39_{\pm 0.44}$	$61.89{\scriptstyle \pm 1.06}$
DER++	$91.14 \pm 0.22$	$92.06 \pm 0.20$	$82.20 \pm 0.89$	$62.67 \pm 1.08$
GPS+DER++	$91.84{\scriptstyle \pm 0.16}$	$92.57_{\pm 0.10}$	$83.53{\scriptstyle \pm 0.64}$	$63.01_{\pm 0.98}$

Table 6: Accuracy of GPS using permutation incorporating DER, DER++ [3] and HAL [4] with small memory sizes  $|\mathcal{M}|$ .

$ \mathcal{M} $	<b>P-MNIST</b> 100	<b>S-CIFAR10</b> 20	<b>S-CIFAR100</b> 200	<b>TinyImageNet</b> 200
HAL	$80.77 \pm 1.31$	$82.56 \pm 2.01$	$52.89 \pm 0.97$	$38.64 \pm 0.89$
GPS+HAL	$81.45{\scriptstyle \pm 0.94}$	$83.74_{\pm 1.75}$	$\boldsymbol{53.57}_{\pm 0.77}$	$39.37_{\pm 0.44}$
DER	$81.72 \pm 1.11$	$85.57 \pm 1.59$	$57.51 \pm 0.60$	$40.21 \pm 0.77$
GPS+DER	$82.04_{\pm 0.97}$	$85.68_{\pm 1.49}$	$57.83_{\pm 0.59}$	$40.54_{\pm 0.54}$
DER++	$83.57{\scriptstyle \pm 0.59}$	$83.45{\scriptstyle \pm 1.76}$	$58.18_{\pm 1.08}$	$40.67 \pm 1.16$
GPS+DER++	$83.86 _{\pm 0.35}$	$83.57_{\pm 1.45}$	$58.15 \pm 0.78$	$40.70{\scriptstyle \pm 1.03}$

### 72 5 New Base Policies based on Curriculum Learning

To further test the power of GPS, we substitute the base policies with two novel memory construction methods designed by us based on curriculum learning [7], ER-CurRes and ER-CurRing-Full. The inspiration of these two methods comes from recent findings that curriculum can help when noisy data are present [11, 8, 10]. We believe data points of future task can be viewed as noisy interference for samples stored in *M*. In these two policies, curricular easy points of each task are picked as candidates for *M*.

#### 79 5.1 Algorithm

**Curricular Easy Samples** We rank data examples from easy to hard based on the implicit cur-80 ricula [11]. Specifically, we first record the learned epochs as an attribute of an example, which is 81 the earliest epoch in training where a model correctly predicts this example for that and subsequent 82 epochs till now. As the learned epoch is a positive integer attribute, it is defined as a subset of the 83 totally ordered set  $\mathbb{Z}^+$ . We also record the current loss of each example as another attribute. The loss 84 attribute is defined as a subset of the totally ordered set  $\mathbb{R}$ . The ranking of examples is based on the 85 lexicographical order on the Cartesian product of the two attributes, *i.e.*, first sorting the examples by 86 the learned epoch attribute, and then ordering examples within the same ranked epoch by their losses. 87 As a result, each training example of a task would be associated with an unique ranking. Also, the 88 ranking would be updated after each epoch.

Algorithm 4 ER-CurRes (for a single task)

**Input:** Reservoir memory buffer  $\mathcal{M}$ ; Number of epochs k; Task distribution P; Dataset size  $|\mathcal{D}|$ ; Batch size B; Portion of easy data  $\gamma$ ; Model parameters  $\theta$ ; Seen examples N. Initialize a random easy pool  $P^{\text{easy}}$  from Pfor  $ep \in \{1, ..., k\}$  do for  $iter \in 1, ..., |\mathcal{D}|/B$  do Sample a batch  $B_P$  from P and a batch  $B_M$  from  $\mathcal{M}$ Update  $\theta$  with  $B_P \cup B_M$ if  $ep \leq \lfloor k/2 \rfloor$  then Update  $\mathcal{M}$  with a probability  $|\mathcal{M}|/N$  for each examples in  $B_P$ N = N + 1else Update  $\mathcal{M}$  with a probability  $|\mathcal{M}|/(\gamma * N)$  for each examples in  $B_P \cap P^{\text{easy}}$  $N = N + 1/\gamma$ end if end for Update  $P^{\text{easy}}$ : order examples based on the implicit curriculum and select the first  $\gamma |\mathcal{D}|$ end for Select the memory for the current task as  $\mathcal{M}_{now}$  and the memory for all previous tasks as  $\mathcal{M}_{past}$ for  $idx \in \mathcal{M}_{now}$  do if  $idx \notin P^{\text{easy}}$  then Replace the slot in  $\mathcal{M}_{now}$  with samples from  $(P^{easy} - \mathcal{M}_{now})$ end if end for **Return** Updated  $\theta$  and  $\mathcal{M} = \mathcal{M}_{now} \cup \mathcal{M}_{past}$ 

89

**ER-CurRes** The ER-CurRes algorithm is shown in Algorithm. 4. Different from ER-Res which stores examples sampled from the whole distribution  $P_i$  of a task  $t_i$  [6], we sample subset of data points from an "easy pool", *i.e.*  $P_i^{\text{easy}}$ . The size is calculated by the dataset size  $|\mathcal{D}_i|$  multiplying a hyperparameter  $\gamma$ , whose value is reported for each evaluation dataset in Section 6. Compared to taking the top few easiest points of each class, sampling from the pool utilizes the benefits of randomness [2, 11]. Suppose we train a total of k epochs of task  $t_i$ , in order to obtain a smooth transition and the samples based on a more stable curriculum ranking, we take the examples obeying ER-Res policy (*i.e.*, samples from  $P_i$ ) for the first  $\lceil k/2 \rceil$  epochs, and take the examples obeying Algorithm 5 ER-CurRing-Full (for a single task)

**Input:** Ring-Full memory buffer  $\mathcal{M}$ ; Number of epochs k; Task distribution P; Dataset size  $|\mathcal{D}|$ ; Batch size B; Portion of easy data  $\gamma$ ; Model parameters  $\theta$ . Initialize a random easy pool  $P^{\text{easy}}$  from PReallocate the memory  $\mathcal{M}_{\text{nost}}$  for all previous tasks, and allocate the memory  $\mathcal{M}_{\text{now}}$  for the current task for  $e \in \{1, ..., k\}$  do for  $iter \in 1, ..., |\mathcal{D}|/B$  do Sample a batch  $B_P$  from P and a batch  $B_M$  from  $\mathcal{M}$ Update  $\theta$  with  $B_P \cup B_M$ if  $e \leq \lfloor k/2 \rfloor$  then Update  $\mathcal{M}_{now}$  with  $B_P$ else Update  $\mathcal{M}_{now}$  with  $B_P \cap P^{easy}$ end if end for Update  $P^{\text{easy}}$ : order examples based on the implicit curriculum and select the first  $\gamma$ portion of each class end for for  $idx \in \mathcal{M}_{now}$  do if  $idx \notin P^{\text{easy}}$  then Replace the slot in  $\mathcal{M}_{now}$  with samples from  $(P^{easy} - \mathcal{M}_{now})$ end if end for **Return:** Updated  $\theta$  and  $\mathcal{M} = \mathcal{M}_{now} \cup \mathcal{M}_{past}$ 

ER-CurRes policy (*i.e.*, samples from  $P_i^{\text{easy}}$ ) for the last  $\lfloor k/2 \rfloor$  epochs. When we finish training, we 98 replace the examples of task  $t_i$  in the memory which are not from  $P_i^{\text{easy}}$ . 99

ER-CurRing-Full Likewise, as shown in Algorithm. 5, ER-CurRing-Full follows the ER-Ring-Full 100 strategy in the first  $\lfloor k/2 \rfloor$  epochs to fill a FIFO memory [6]. After that, we substitute the memory 101 slots with points from the easy pool of each observed class. In the construction of the easy pool, 102 instead of taking the easiest  $\gamma |\mathcal{D}|$  examples as in ER-CurRes, we use the easiest  $\gamma |\mathcal{D}|/C$  examples of 103 each class based on the ranking, where C is the number of classes. 104

#### 5.2 Experimental Results of GPS w/ Cur 105

The accuracy comparison between GPS w/ Cur and GPS w/ ER-CurRes, ER-CurRing-Full are shown 106 in Table 7. From the table, we can see GPS w/ Cur outperforms both ER-CurRes and ER-CurRing-107 Full in both datasets. Moreover, it shows leveraging curriculum in base policies of GPS can further 108

improve its performance compared to its plain version. 109

$ \mathcal{M} $	<b>P-MNIST</b> 1000	<b>S-CIFAR10</b> 200	<b>S-CIFAR100</b> 2000	<b>TinyImageNet</b> 2000
ER-CurRes ER-CurRing-Full GPS w/ Cur	$\begin{array}{c} 86.78 {\scriptstyle \pm 0.49} \\ 86.16 {\scriptstyle \pm 0.49} \\ \textbf{88.35} {\scriptstyle \pm 0.18} \end{array}$	$\begin{array}{c} 92.47 {\scriptstyle \pm 0.20} \\ 91.70 {\scriptstyle \pm 0.50} \\ \textbf{93.58} {\scriptstyle \pm 0.17} \end{array}$	$\begin{array}{c} 81.38 {\scriptstyle \pm 0.51} \\ 81.16 {\scriptstyle \pm 0.65} \\ \textbf{82.34} {\scriptstyle \pm 0.79} \end{array}$	$\begin{array}{c} 61.89{\scriptstyle\pm0.03} \\ 61.03{\scriptstyle\pm0.42} \\ \textbf{62.88}{\scriptstyle\pm0.03} \end{array}$
$ \mathcal{M} $	<b>P-MNIST</b> 100	<b>S-CIFAR10</b> 20	<b>S-CIFAR100</b> 200	<b>TinyImageNet</b> 200
ER-CurRes ER-CurRing-Full <b>GPS w/ Cur</b>	$\begin{array}{c} 65.34 {\scriptstyle \pm 0.69} \\ 66.42 {\scriptstyle \pm 1.03} \\ \textbf{66.92 {\scriptstyle \pm 0.54}} \end{array}$	$\begin{array}{c} 81.59{\scriptstyle\pm3.23}\\ 81.54{\scriptstyle\pm2.46}\\ \textbf{81.68}{\scriptstyle\pm1.98}\end{array}$	$\begin{array}{c} 65.43{\scriptstyle\pm1.37}\\ 66.73{\scriptstyle\pm0.12}\\ \textbf{67.08}{\scriptstyle\pm0.36}\end{array}$	$\begin{array}{c} 40.75 \scriptstyle{\pm 0.98} \\ 41.54 \scriptstyle{\pm 0.57} \\ \textbf{41.80} \scriptstyle{\pm 0.49} \end{array}$

Table 7: Accuracy of GPS using curriculum-based policies vs. the corresponding baselines.

Dataset	Method	Parameters
P-MNIST	CurER-Res CurER-Ring-Full GPS w/ Cur DER DER++ GPS+DER GPS+DER++ HAL GPS+HAL oEWC GSS OGD	$\begin{array}{c} \gamma: \ 0.2 \\ \gamma: \ 0.1 \\ \gamma: \ 0.2 \\ \alpha: \ 0.5 \\ \alpha: \ 1.0 \ \beta: \ 0.5 \\ \alpha: \ 0.5 \\ \alpha: \ 1.0 \ \beta: \ 0.5 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \ \gamma: \ 0.1 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \ \gamma: \ 0.1 \\ \lambda: \ 0.7 \ \gamma: \ 1.0 \\ gmbs: \ 10 \ nb: \ 1 \\ stored \ gradients: \ 100/task \ (perm) \end{array}$
S-CIFAR-10	CurER-Res CurER-Ring-Full GPS w/ Cur DER DER++ GPS+DER GPS+DER++ HAL GPS+HAL oEWC iCaRL GSS	$\begin{array}{c} \gamma: \ 0.2 \\ \gamma: \ 0.1 \\ \gamma: \ 0.2 \\ \alpha: \ 0.3 \\ \alpha: \ 0.1 \ \beta: \ 0.5 \\ \alpha: \ 0.3 \\ \alpha: \ 0.1 \ \beta: \ 0.5 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \ \gamma: \ 0.1 \\ \lambda: \ 0.1 \ \beta: \ 0.5 \ \gamma: \ 0.1 \\ \lambda: \ 0.7 \ \gamma: \ 1.0 \\ wd: \ 0 \\ gmbs: \ 32 \ nb: \ 1 \end{array}$
S-CIFAR-100	CurER-Res CurER-Ring-Full GPS w/ Cur DER DER++ GPS+DER GPS+DER HAL GPS+HAL oEWC iCaRL GSS	$\begin{array}{c} \gamma: \ 0.2 \\ \gamma: \ 0.1 \\ \gamma: \ 0.2 \\ \alpha: \ 0.5 \\ \alpha: \ 0.5 \\ \beta: \ 0.5 \\ \alpha: \ 0.5 \\ \beta: \ 0.5 \\ \alpha: \ 0.5 \\ \beta: \ 0.5 \\ \lambda: \ 0.1 \\ \beta: \ 0.5 \\ \gamma: \ 0.1 \\ \lambda: \ 0.1 \\ \beta: \ 0.5 \\ \gamma: \ 0.1 \\ \lambda: \ 0.7 \\ \gamma: \ 1.0 \\ wd: \ 10^{-5} \\ gmbs: \ 32 \ nb: \ 1 \end{array}$
TinyImageNet	CurER-Res GPS w/ Cur CurER-Ring-Full DER DER++ GPS+DER GPS+DER++ HAL GPS+HAL oEWC iCaRL	$\begin{array}{l} \gamma: \ 0.2 \\ \gamma: \ 0.2 \\ \gamma: \ 0.1 \\ \alpha: \ 0.1 \\ \alpha: \ 0.1 \\ \beta: \ 0.5 \\ \alpha: \ 0.1 \\ \alpha: \ 0.1 \\ \beta: \ 0.5 \\ \lambda: \ 0.1 \\ \beta: \ 0.5 \\ \gamma: \ 0.1 \\ \lambda: \ 0.1 \\ \beta: \ 0.5 \\ \gamma: \ 0.1 \\ \lambda: \ 0.7 \\ \gamma: \ 1.0 \\ wd: \ 10^{-5} \end{array}$

Table 8: Other hyperparameters used in our experiments.

# **110 6 Experimental Details**

# 111 6.1 Simulation Details

In experiments, we set the number of examples in each synthesized pseudo-task the same as the size of the memory buffer, *i.e.*, if  $|\mathcal{M}| = 1000$ , we then generate 1000 examples for each pseudo-task. For computational efficiency, we set the number of training epochs small in the simulated training process. We train 1 epoch for pseudo-tasks synthesized in the P-MNIST dataset, 3 epochs for pseudo-tasks in S-CIFAR-10, S-CIFAR-100 and TinyImageNet. As for the batch size, the optimizer and the learning the simulation process, they are all the same as in the real training process.

#### 118 6.2 Other Hyperparameters

We disclose the experimental hyperparameters values not reported in the main manuscript in Table 8. In the table,  $\gamma$  in the 'Cur-' series methods is the easy pool ratio of the curriculum-based policies as we discussed in Section 5, while other symbols refer to the respective methods. In all the experimental evaluation by accuracy, reported numbers are averaged over 5 runs.

## 123 6.3 Time Measurement

We measure our training and simulation time for each dataset in a single NVIDIA Tesla K80 GPU for fair comparison. The time we report is the total processing time averaged on 5 runs, assessed in wall-clock time (seconds) at the end of the last task and then converted into minutes.

# 127 **References**

- [1] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection
  for online continual learning, 2019.
- 130 [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning.
- [3] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark
  experience for general continual learning: a strong, simple baseline, 2020.
- [4] Arslan Chaudhry, Albert Gordo, Puneet K. Dokania, Philip Torr, and David Lopez-Paz. Using
  hindsight to anchor past knowledge in continual learning, 2021.
- [5] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient
  lifelong learning with a-gem, 2019.
- [6] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K.
  Dokania, Philip H. S. Torr, and Marc'Aurelio Ranzato. On tiny episodic memories in continual learning, 2019.
- [7] Jeffrey Elman. Learning and development in neural networks: the importance of starting small.
  *Cognition*, 48:71–99, 08 1993.
- [8] Melike Nur Mermer and Mehmet Fatih Amasyali. Training with growing sets: A simple alternative to curriculum learning and self paced learning, 2018.
- [9] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl:
  Incremental classifier and representation learning, 2017.
- [10] Sebastian Ruder and Barbara Plank. Learning to select data for transfer learning with bayesian
  optimization, 2017.
- <sup>148</sup> [11] Xiaoxia Wu, Ethan Dyer, and Behnam Neyshabur. When do curricula work?, 2021.